**Names: Nehil Joshi, Harshil Shah, Jashesh Mehta, Hasti Panchal.**
**Group 36**
**Subject: CSP 571 Data Preparation & Analysis**

# Final Project: Data Scientist Salary Analysis

## Analyzing Salary Trends and Building Predictive Models

## Introduction

- Background:
  In today's competitive job market, salary data for roles such as Data Scientists is invaluable for both job seekers and companies aiming to attract the best talent. This project explores trends and patterns in Data Scientist salaries across different industries, locations, and job titles.

- Objectives:
  The objective of this project is to:

    o Analyze the salary data for Data Scientist roles.

    o Visualize salary trends and correlations between key features.

    o Build predictive models to estimate salaries.

- Dataset Overview:
  The dataset used in this project contains information from 742 job postings for Data Scientist roles. Key features include job titles, company details, skills, and salary ranges.

## Data Cleaning and Preprocessing

- Handling Missing Data:
  Missing values in the dataset were handled by replacing invalid values, such as -1, with NaN. The Avg Salary(K) feature was created by averaging the lower and upper salary bounds.

- Exploratory Data Analysis: To understand the dataset better, basic statistics and missing values were explored using df.describe(), df.info(), and df.isnull().sum().
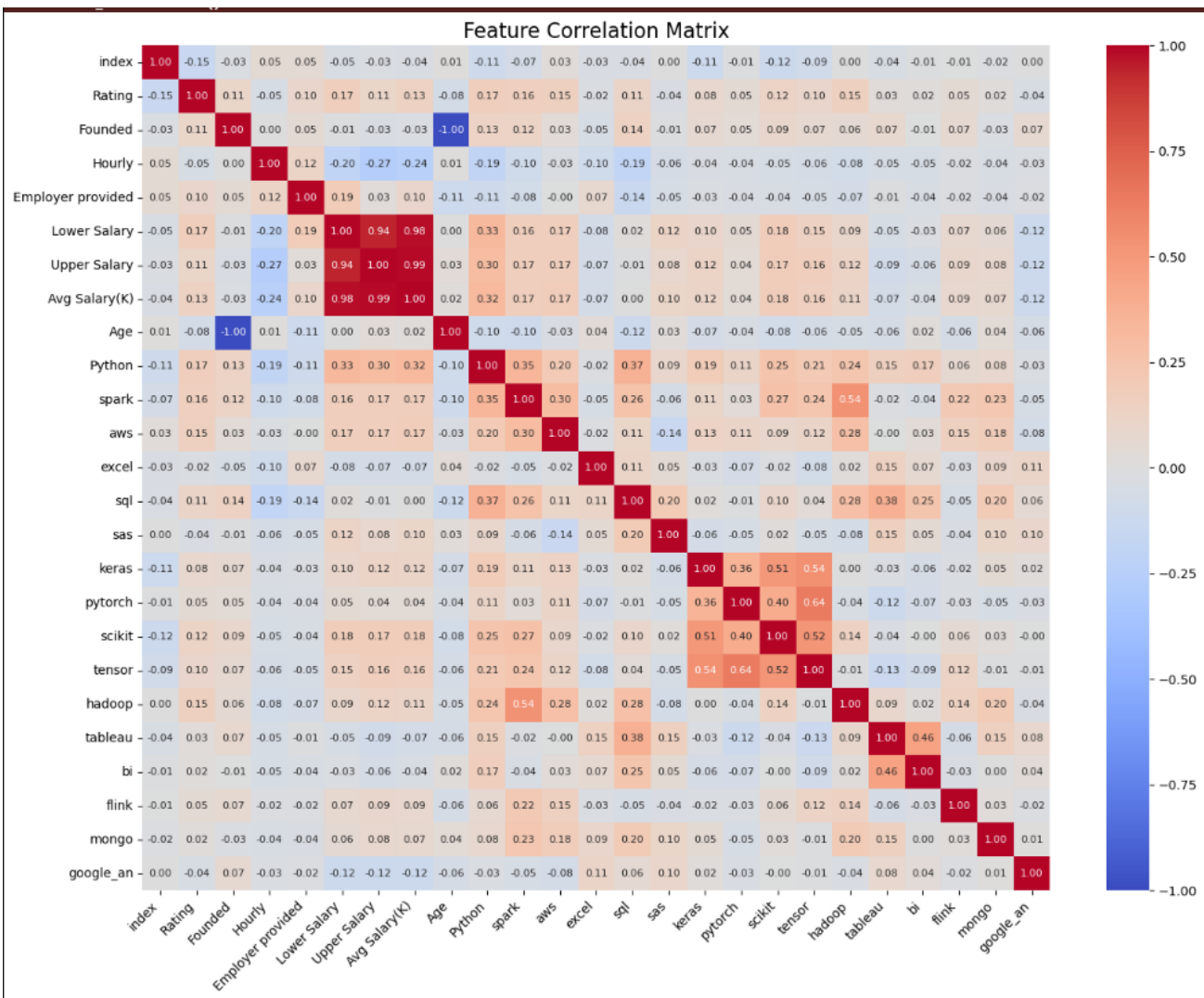
- Code:

```
# Replace -1 values with NaN
df['Rating'].replace(-1, np.nan, inplace=True)
df['Founded'].replace(-1, np.nan, inplace=True)
# Create Avg Salary column
df['Avg Salary(K)'] = (df['Lower Salary'] + df['Upper Salary']) / 2
```

- Correlation Analysis:

  A correlation matrix was calculated to identify relationships between numerical variables, especially between Avg Salary(K) and other features such as Lower Salary and Upper Salary.
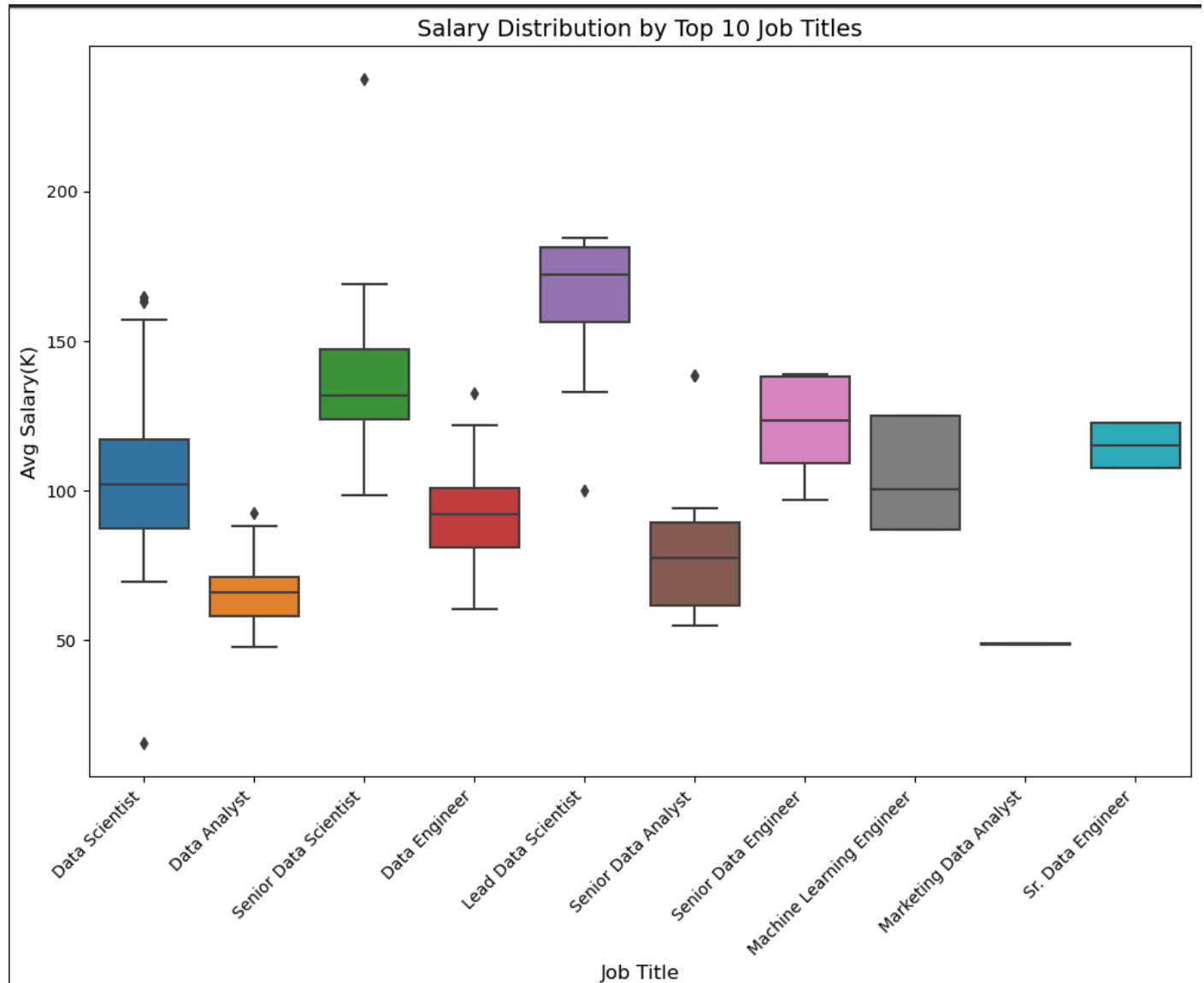
  - **Graph: Correlation Heatmap**



Feature Correlation Matrix

# Exploratory Data Analysis (EDA)

- Salary Distribution by Job Titles:
  The distribution of salaries by job title was explored using a box plot to highlight the variation in salaries for different roles. Job titles like Senior Data Scientist and Machine Learning Engineer exhibited higher median salaries compared to roles like Data Analyst.
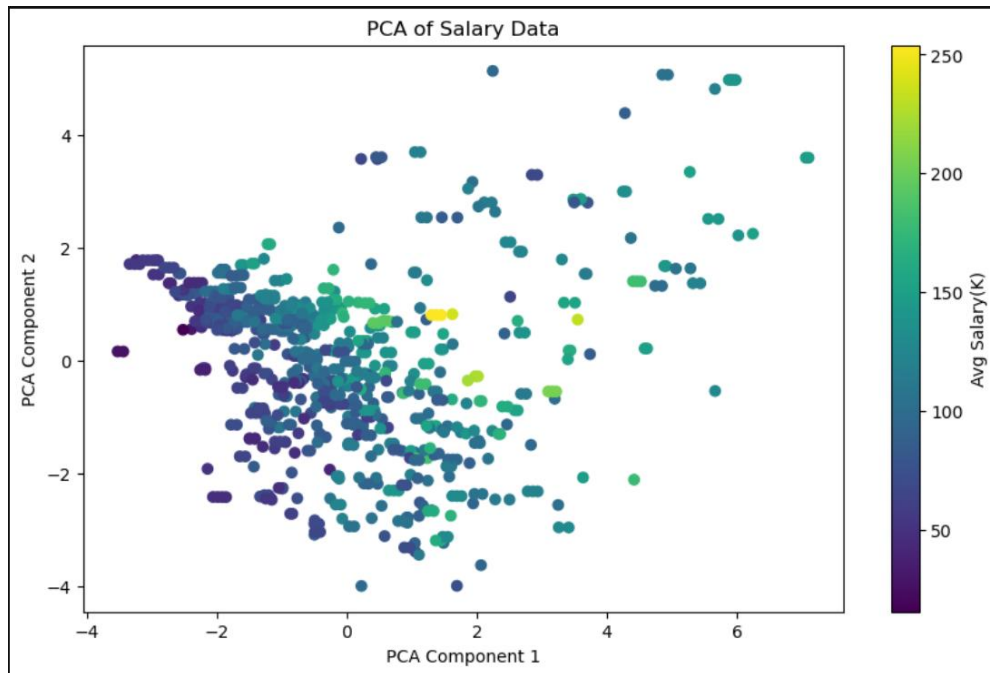
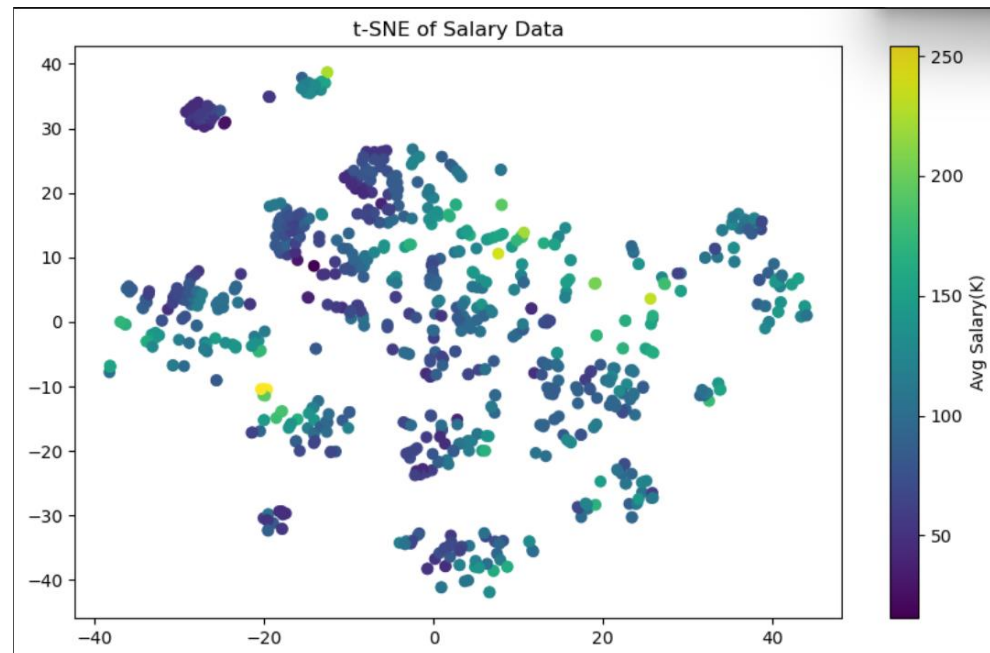  - **Graph: Salary Distribution by Top 10 Job Titles**

- Dimensionality Reduction:

  PCA and t-SNE were applied to reduce the high-dimensional feature space into two dimensions to visualize the data's underlying structure.

  o **Graph: PCA of Salary Data**

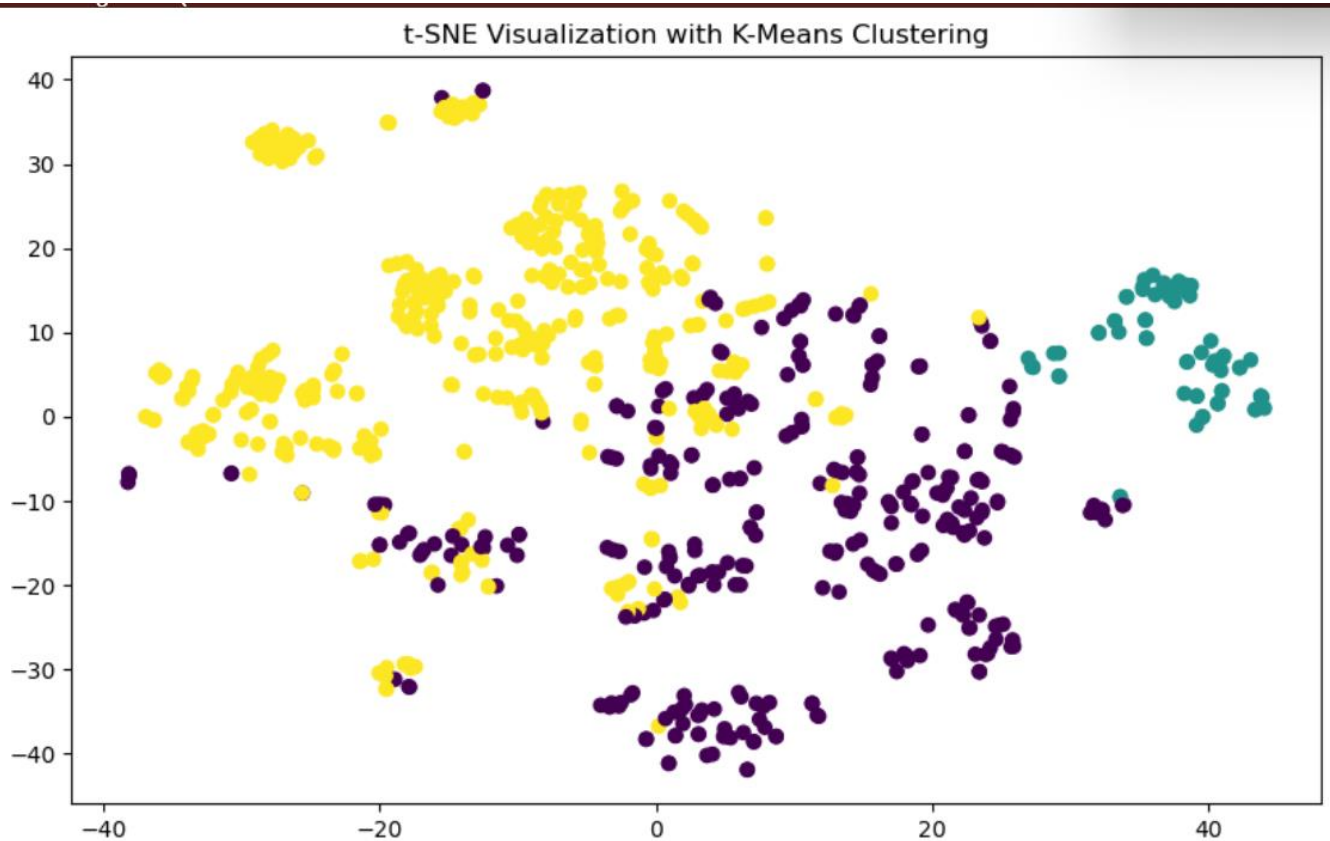

  o **Graph: t-SNE of Salary Data**

# Clustering and Predictive Modeling

- ## Clustering:
  K-means clustering was applied to segment the dataset into three clusters based on the features. The t-SNE visualization with clusters helps understand how data points are grouped.

  - ### Graph: t-SNE with K-Means Clusters



- ## Predictive Modeling:
  Several machine learning models were trained to predict salaries, including:

  - ### Linear Regression

  - ### Lasso and Ridge Regression

  - ### Gradient Boosting Regressor

After training, the models were evaluated using metrics like $R^2$, Mean Absolute Error (MAE), and Mean Squared Error (MSE).

- o **Code for Model Training (Example: Linear Regression):**
  ```
  model = LinearRegression()
  model.fit(X_train_selected, y_train)
  y_pred = model.predict(X_test_selected)
  print("MAE:", mean_absolute_error(y_test, y_pred))
  print("MSE:", mean_squared_error(y_test, y_pred))
  print("R2 Score:", r2_score(y_test, y_pred))
  ```

- Model Results:
  - o Linear Regression achieved an $R^2$ of 0.82.
  - o Lasso and Ridge regularization models were used to avoid overfitting, with $R^2$ scores of 0.83 and 0.81, respectively.
  - o Gradient Boosting outperformed all, with an $R^2$ of 0.88, making it the best model for salary prediction.

# Insights and Conclusions

- Key Insights:
  - o Specialized roles such as Senior Data Scientist and Machine Learning Engineer tend to offer higher salaries.
  - o Data from large companies, particularly in tech and finance, indicates higher salary ranges.
  - o Skills such as Python, SQL, and cloud technologies like AWS are highly correlated with higher salaries.

- Recommendations for Job Seekers:
  - o Job seekers should focus on developing in-demand technical skills.
  - o Targeting specialized roles like Senior Data Scientist or Machine Learning Engineer can significantly increase earning potential.

- Conclusion: This project successfully explored and analyzed Data Scientist salaries. By visualizing key trends and building predictive models, it provides insights that can assist both job seekers and employers in understanding salary dynamics. The Gradient Boosting model proved to be the most effective in predicting salaries, offering a robust tool for future salary estimations.